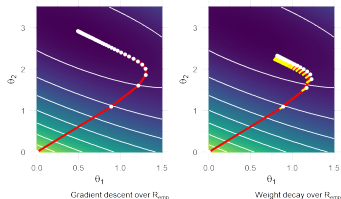


Weight Decay and L2



- Understand why $L2$ regularization in combination with gradient descent is equivalent to weight decay
- Understand how weight decay changes the optimization trajectory

WEIGHT DECAY VS. L2 REGULARIZATION

Let us optimize the L_2 -regularized risk of a model $f(\mathbf{x} \mid \theta)$

$$\min_{\theta} \mathcal{R}_{\text{reg}}(\theta) = \min_{\theta} \mathcal{R}_{\text{emp}}(\theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

by gradient descent. The gradient is

$$\nabla_{\theta} \mathcal{R}_{\text{reg}}(\theta) = \nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta) + \lambda \theta.$$

We iteratively update θ by step size α times the negative gradient

$$\begin{aligned} \theta^{[\text{new}]} &= \theta^{[\text{old}]} - \alpha \left(\nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta^{[\text{old}]}) + \lambda \theta^{[\text{old}]} \right) \\ &= \theta^{[\text{old}]} (1 - \alpha \lambda) - \alpha \nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta^{[\text{old}]}). \end{aligned}$$

The term $\lambda \theta^{[\text{old}]}$ causes the parameter (**weight**) to **decay** in proportion to its size. This is a very well-known technique in deep learning - and simply L_2 regularization in disguise (for gradient descent).



WEIGHT DECAY VS. L2 REGULARIZATION / 2

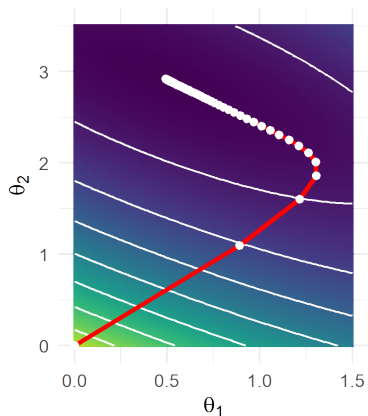
Caveat: Equivalence of weight decay and $L2$ only holds for (S)GD!

- [Hanson and Pratt 1988](#) originally define WD “decoupled” from gradient-updates $\alpha \nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta^{[\text{old}]})$ as
$$\theta^{[\text{new}]} = \theta^{[\text{old}]}(1 - \lambda') - \alpha \nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta^{[\text{old}]})$$
- This is equivalent to modern WD/ $L2$ (last slide) using reparameterization $\lambda' = \alpha \lambda$
- Consequence: if there is optimal λ' , then optimal $L2$ penalty is tightly coupled to α as $\lambda = \lambda' / \alpha$ (and vice versa)
- [Loshchilov and Hutter 2019](#) show no equivalence of $L2$ and WD possible for adaptive methods like Adam (Prop. 2)
- In many cases where SGD+ $L2$ works well, Adam+ $L2$ underperforms due to non-equivalence with WD
- They propose a variant of Adam decoupling WD from gradient updates (AdamW), increasing performance over Adam+ $L2$

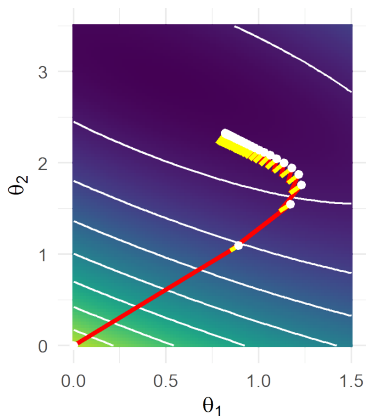


WEIGHT DECAY VS. L2 REGULARIZATION / 3

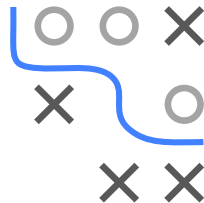
When we use weight decay, we follow the steepest slope of \mathcal{R}_{emp} as for gradient descent, but in every step, we are pulled back to the origin.



Gradient descent over \mathcal{R}_{emp}

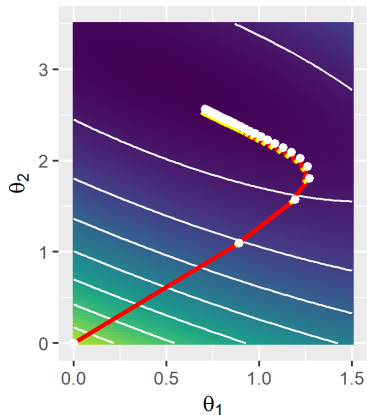


Weight decay over \mathcal{R}_{emp}

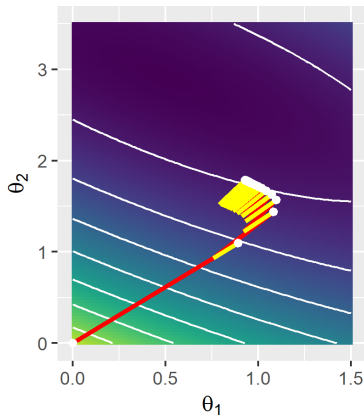


WEIGHT DECAY VS. L2 REGULARIZATION / 4

How strongly we are pulled back to the origin (for a fixed stepsize α) depends only on λ as long as the procedure converges:



Weight decay (small λ) over R_{emp}



Weight decay (large λ) over R_{emp}

